

Using shortest path to discover criminal community

Pritheega Magalingam^{a,b,*}, Stephen Davis^a, Asha Rao^a

^a*School of Mathematical and Geospatial Sciences, RMIT University, Melbourne, Australia, GPO Box 2476, Melbourne, Victoria 3001.*

^b*Advanced Informatics School, Level 5, Menara Razak, Universiti Teknologi Malaysia, Jalan Semarak, 54100 Kuala Lumpur, Malaysia.*

Abstract

Extracting communities using existing community detection algorithms yields dense sub-networks that are difficult to analyse. Extracting a smaller sample that embodies the relationships of a list of suspects is an important part of the beginning of an investigation. In this paper, we present the efficacy of our shortest paths network search algorithm (SPNSA) that begins with an ‘algorithm feed’, a small subset of nodes of particular interest, and builds an investigative sub-network. The algorithm feed may consist of known criminals or suspects, or persons of influence. This sets our approach apart from existing community detection algorithms. We apply the SPNSA on the Enron Dataset of e-mail communications starting with those convicted of money laundering in relation to the collapse of Enron as the algorithm feed. The algorithm produces sparse and small sub-networks that could feasibly identify a list of persons and relationships to be further investigated. In contrast, we show that identifying sub-networks of interest using either community detection algorithms or a k-Neighbourhood approach produces sub-networks of much larger size and complexity. When the 18 top managers of Enron were used as the algorithm feed, the resulting sub-network identified 4 convicted criminals that were not managers and so not part of the algorithm feed. We also directly tested the SPNSA by removing one of the convicted criminals from the algorithm feed and re-running the algorithm; in 5 out of 9 cases the left out criminal occurred in the resulting sub-network.

Keywords: Criminal network, Shortest path, Leave-one-out, Trust, Suspect, Investigation

1. Introduction

Retrieving a criminal network from an organised crime incident is an important part of crime investigation. This task is a difficult one, mainly because of the involvement of a variety of criminals who play myriad roles (Basu, 2014; Didimo et al., 2011). In addition to drug trafficking and money laundering, organised crime includes hijacking and equipment smuggling. The task of the criminal investigator is further

*Corresponding author. Tel.: +61 3 9925 1843

Email addresses: pritheega.magalingam@rmit.edu.au (Pritheega Magalingam), asha@rmit.edu.au (Asha Rao)

hampered by the mass of data needing to be searched with an important part of the start of an investigation being the identification of a smaller sample that embodies the relationships within the criminal participants. In (Magalingam et al., 2014), we presented an algorithm designed to extract a smaller, more manageable, network of possible relationships from a large dataset of interactions. In this paper, we further develop this algorithm and show that it performs well in a variety of scenarios, and is able to extract meaningful sub-networks for a criminal investigator to start an investigation. We show that this algorithm performs better than known community detection algorithms (Pons, 2006; Clauset et al., 2004; Newman, 2006), as well as k-neighbourhood detection methods (Zhou and Pei, 2011).

In the past, extracting criminal associations from raw data has required preliminary information of such relationships, while building a network from such, known, relationships has been done manually (Basu, 2014; Didimo et al., 2011; Christin et al., 2010; Oatley and Crick, 2014). For example, Nadji et al. (2013) produce a network of known fraudulent infrastructure by creating links between IP addresses using known attack signatures garnered from passive domain name server and several other sources for malicious activities. Krebs (2002) builds edges between known hijackers of the 9-11 terrorist attacks by manually gathering data from online news articles. The edges, or links, are created based of information such as whether the two persons went to the same school, grew up in the same locality, etc. Oatley and Crick (2014) follow a similar track, using associations such as partner, sibling, cohabitant, to build a relationship network among the members of different UK crime gangs. Clearly, the above methods are time-consuming, and a faster, more automated process of building a relationship network would be very useful for investigators of criminal activities.

We present such an algorithm, which can be run on a large dataset of interactions, to build a more practicable sub-network of known criminals suitable for further investigation. We use the publicly available Enron Dataset (Cohen, 2009), which contains all email communications before and after the collapse of this large company in 2001. This dataset is appropriate for this exercise, as ten people connected with Enron were subsequently convicted of money laundering (Securities and Release, 2004).

The structure of the rest of the paper is as follows: In the next section, we describe the Enron dataset in more detail, give the process by which we start the isolation of specific email groupings, compare the connections between the ten criminals in two different email sub-networks, and describe our algorithm. Section 3 gives the results of applying existing community detection algorithms as well as the k-nearest neighbour method, to the Enron dataset to identify the community that the criminals belong. In the section after this, we apply our shortest paths network search algorithm to the two email sub-networks previously identified and compare the results to those obtained by applying the existing community detection algorithms. The penultimate section details the application of our algorithm to the different scenarios that an investigator may encounter. Finally we give the conclusion.

2. Background

This section describes the preliminary analysis of the Enron email dataset, the people who were convicted of money laundering crime, the identification of criminal communication links and the criminal sub-network formation methods.

2.1. Preliminary analysis of dataset

The Enron email dataset contains 1,887,305 email transactions (Cohen, 2009) that were sent using the fields ‘TO’, ‘CC’ or ‘BCC’. Out of these emails, 16,116 are senders of the emails and 68,203 are receivers of the emails. The Enron email dataset contains a mix of internal and external email transactions. Within the 16,116 email senders, 5,831 email transactions are from email addresses that are Enron company email accounts having the name ‘enron’ in their email address and the rest of the addresses are external, for example *andrew.fastow@ljminvestments.com*, *anitatr@earthlink.net*, etc. In order to process this large number of emails, we start by extracting the emails sent and received in the last 8 years of Enron - from 1995 to 2002 (Salter, 2008). We clean the data by removing the irrelevant email transactions such as email addresses that have numbers and characters for example ‘5673@aol.com’, that end with airline company name for example ‘@aircanada.com’, that end with ‘xpedia.com’, ‘amazon.com’ and other auto response emails.

Several prior works propose ways of extracting criminal networks in the form of associations between texts or people (Basu, 2014; Krebs, 2002). Mining relevant terms from a large volume of police incident summaries and assigning the co-occurrence frequency as a weight to each term is used by (Chen et al., 2004) to design a criminal network while Yang and Ng (2007) use web crawlers to gather identities associated with certain crime related topics in web blog pages and represent them as a network. Similarly, in order to identify criminal cliques, Iqbal et al. (2012) perform chat topic analysis and certain entities that belong to the same chat session are formed into a clique. Louis and Engelbrecht (2011) conduct text mining on passages of a mystery novel to show the association between words, in the form a graph, leading to the identification of murders. In (Anwar and Abulaish, 2012), posts that promote hate and violence in certain dark web forums are grouped in different cliques using an algorithm that measures similarity based on content, time, author and title.

Using keywords as a tool for isolating criminal networks is a problem especially when electronic documents, chat messages, web blogs or emails contain incomplete information or could mislead detection algorithms (Murynets and Piqueras Jover, 2012; Keila and Skillicorn, 2005). Consequently, we choose to ignore the content of the various emails being exchanged between the criminals and propose a very different way to start the building of a criminal network, by considering the type of emails based on recipient fields. As detailed in (Magalingam et al., 2014), we separate the emails with at least one BCC recipient because the existence of a bcc in an email, could indicate a trust relationship (Fox and Schaefer, 2012). While ‘to’

and ‘cc’ recipients are visible to all recipients, as (McDowell and Householder, 2009) and (Bogawar and Bhoyar, 2012) point out, there is something inherently secretive about adding a ‘bcc’ recipient to an email. We explore the suspect’s secret connections in the group of emails that consists of all those with one or more bcc-ed recipients and compare our result with the connections found using email transactions that have recipients only in the ‘TO’ and ‘CC’ fields.

The emails are divided into two groups. The first group is made up of email transactions that have recipients in the ‘TO’ and ‘CC’ fields. These emails do not contain any BCC recipients. The network formed using the ‘TO’ and ‘CC’ email transactions has 26,027 nodes and 1,048,572 edges with an average degree of 80.58. Henceforth, we refer to this network as Netgraph and each node ID is now called Net. ID. The relationship between the nodes’ degrees and their frequencies is displayed in the log-log plot of the degree distribution (Figure 1 (a)) .

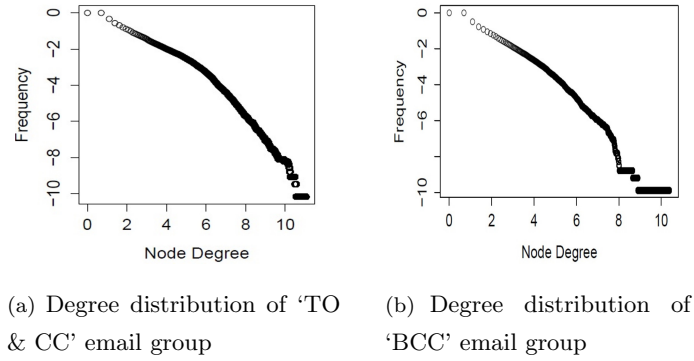


Figure 1: The figures show that the degree distributions for both Netgraph (a) and BCC Netgraph (b) are heavy-tailed with many nodes being of low degree and some nodes being highly connected.

The second group of emails consists of all those with one or more bcc-ed recipients. The network formed using the BCC email transactions is called the BCC Netgraph with each node in this BCC Netgraph given a BCCNet. ID. The BCC Netgraph contains 19,716 nodes and 238,761 edges. The BCC Netgraph has an average degree of 24.22 and the log-log plot of the degree distribution is shown in Figure 1 (b). As is evident from Figure 1 (a) and (b), both networks have many nodes with low degree and a few nodes with very high degree indicating the degree distributions are heavy tailed. Sub-networks are then constructed with each of these large networks using our shortest paths network search algorithm (see Section 2.4). Next the persons convicted of money laundering in relation to the collapse of Enron are listed along with their Net. IDs and BCCNet. IDs.

2.2. Enron money laundering criminals

Ten people were convicted of money laundering in relation to the collapse of Enron (Brickey, 2003). Table 1 below shows the Net. ID of the criminals appearing within Netgraph as well as the BCCNet. ID of

those appearing within the BCC Netgraph.

Table 1: Enron money laundering criminals

Name	Net. ID	BCCNet. ID	Email Address
Andrew Fastow	1472	686	andrew.fastow@enron.com
Andrew Fastow	-	687	andrew.fastow@ljminvestments.com
Lea Fastow	17589	11010	lfastow@pop.pdq.net
Lea Fastow	17588	11009	lfastow@pdq.net
Kevin Hannon	16202	10068	kevin.hannon@enron.com
Kenneth Rice	16115	9994	kenneth.rice@enron.com
Rex Shelby	23983	15224	rex.shelby@enron.com
Rex Shelby	23985	15225	rex_shelby@enron.net
A. Khan	-	205	adnankkhan@hotmail.com
Michael Kopper	20217	12708	michael.kopper@enron.com
Ben Glisan	-	1369	ben.glisan@enron.com
Joe Hirko	14052	8716	joe.hirko@enron.com
S. Yaeger	-	861	anne.yaeger@enron.com

Table 1 gives the list of e-mail accounts associated with criminals involved in the Enron money laundering crime (Brickey, 2003; Cohen, 2009). The IDs in the table are computer generated numbers assigned to distinct email addresses based on the type of network. The Net. ID refers to the email addresses of criminals in the Netgraph while the BCCNet. ID refers to the email addressess of criminals in the BCC email network.

The multiple email addresses of the criminals, leading to multiple, different IDs, are preserved as some of these criminals, for example Andrew Fastow (BCCNet. ID 687), A. Khan (BCCNet. ID 205), Ben Glisan (BCCNet. ID 1369) and S. Yaeger (BCCNet. ID 861), did not occur in the Netgraph but were present in the BCC Netgraph. Each of these email addresses occur in distinct email transactions. We now identify the length of the shortest paths between these criminals in both the Netgraph and the BCC Netgraph.

2.3. Distribution of criminal links in Netgraph and BCC Netgraph

Here, an analysis is conducted to compare the connections formed in the Netgraph and BCC Netgraph. Network measures are often used to quantify network structures, for example, the number of vertices in a network measures the size of the network, a vertex's degree could be used to show whether it is strong or weak, the shortest path length measures the distance between vertices, centrality measures demonstrate the level of importance of vertices, etc. (Newman, 2010). We use the length of the shortest paths and the

degree of the node. Both the Netgraph and the BCC Netgraph are directed graphs. In order to analyse a criminal's communication links within these graphs, we first calculate the length of the shortest paths from one criminal to another. Tables 2 and 3 shows the length of shortest paths from criminal to criminal in the directed Netgraph and the BCC Netgraph respectively. The directed Netgraph has an average path length 3.203258 while the directed BCC Netgraph's average path length is 4.445533.

Table 2: Shortest path length from criminal to criminal in the directed Netgraph

Directed Netgraph									
	1472	17589	17588	16202	16115	23983	23985	20217	14052
1472	0	3	3	4	3	2	3	3	3
17589	Inf	0	Inf	Inf	Inf	Inf	Inf	Inf	Inf
17588	Inf	Inf	0	Inf	Inf	Inf	Inf	Inf	Inf
16202	Inf	Inf	Inf	0	Inf	Inf	Inf	Inf	Inf
16115	Inf	Inf	Inf	Inf	0	Inf	Inf	Inf	Inf
23983	2	3	3	4	2	0	2	2	3
23985	Inf	Inf	Inf	Inf	Inf	Inf	0	Inf	Inf
20217	Inf	Inf	Inf	Inf	Inf	Inf	Inf	0	Inf
14052	Inf	Inf	Inf	Inf	Inf	Inf	0	Inf	0

Table 2 shows the lengths of shortest paths from criminal to criminal in the directed Netgraph.

Table 3: Shortest path length from criminal to criminal in the directed BCC Netgraph

Directed BCC Netgraph													
	686	687	11010	11009	10068	12708	1369	9994	8716	15224	15225	861	205
686	0	4	3	3	2	3	2	3	3	2	4	4	Inf
687	Inf	0	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf
11010	3	1	0	3	3	3	3	3	3	4	5	5	Inf
11009	Inf	Inf	Inf	0	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf
10068	2	4	3	3	0	3	3	2	2	3	5	5	Inf
12708	Inf	Inf	Inf	Inf	Inf	0	Inf	Inf	Inf	Inf	Inf	Inf	Inf
1369	1	4	3	3	2	1	0	3	3	3	3	2	Inf
9994	Inf	Inf	Inf	Inf	Inf	Inf	Inf	0	Inf	Inf	Inf	Inf	Inf
8716	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	0	Inf	Inf	Inf	Inf
15224	3	4	3	3	3	3	3	3	3	0	4	4	Inf
15225	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	0	Inf	Inf
861	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	0	Inf
205	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	0

Table 3 shows the lengths of shortest paths from criminal to criminal in the directed BCC Netgraph.

From Tables 2 and 3 it is clear that, in the directed graphs under consideration, only a few criminals have directed paths connecting them to other criminals. If we make the broad assumption that an email sent from A to B implies an *undirected* relationship between A and B then the graphs become undirected. In this case, in BCC Netgraph, 12 of the 13 accounts associated with criminals belong to the same connected component and a path can be found from one criminal's account to another's (see Table 5). The exception is the account associated with A. Khan (adnankkhan@hotmail.com) which belongs to a separate component. The assumption regarding reciprocal relationship seems most appropriate for the trust network (BCC Netgraph) where if A includes B as a BCC recipient there is a personal trust relationship implied between A and B that we assume is reciprocated to some degree.

The shortest path lengths between the criminals in the undirected graphs are shown in Tables 4 and 5. The average path length values of the undirected Netgraph and the undirected BCC Netgraph are 3.264676 and 5.033507 respectively. The average path length of the criminals in the undirected Netgraph and the undirected BCC Netgraph are 2.93 and 3.65 respectively; lower than the average path length of the entire graphs.

Table 4: Shortest path length from criminal to criminal in the undirected Netgraph

Undirected Netgraph									
	1472	17589	17588	16202	16115	23983	23985	20217	14052
1472	0	3	2	3	2	2	3	2	2
17589	3	0	2	4	3	3	3	3	3
17588	2	2	0	4	3	3	4	3	3
16202	3	4	4	0	3	3	4	4	4
16115	2	3	3	3	0	2	3	2	2
23983	2	3	3	3	2	0	3	2	2
23985	3	3	4	4	3	3	0	4	4
20217	2	3	3	4	2	2	4	0	3
14052	2	3	3	4	2	2	4	3	0

Table 4 above shows the lengths of shortest paths from criminal to criminal in undirected Netgraph.

Table 5: Shortest path length from criminal to criminal in the undirected BCC Netgraph

Undirected BCC Netgraph													
	686	687	11010	11009	10068	12708	1369	9994	8716	15224	15225	861	205
686	0	4	3	3	2	2	1	2	3	3	5	4	Inf
687	4	0	1	3	4	5	4	5	5	5	6	7	Inf
11010	3	1	0	2	3	4	3	4	4	4	5	6	Inf
11009	3	3	2	0	3	4	3	4	4	3	5	6	Inf
10068	2	4	3	3	0	2	2	2	2	2	5	5	Inf
12708	2	5	4	4	2	0	1	3	3	3	5	4	Inf
1369	1	4	3	3	2	1	0	2	3	3	4	3	Inf
9994	2	5	4	4	2	3	2	0	2	3	5	5	Inf
8716	3	5	4	4	2	3	3	2	0	2	5	6	Inf
15224	3	5	4	3	2	3	3	3	2	0	4	5	Inf
15225	5	6	5	5	5	5	4	5	5	4	0	6	Inf
861	4	7	6	6	5	4	3	5	6	5	6	0	Inf
205	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	Inf	0

Table 5 above shows the lengths of shortest paths from criminal to criminal in undirected BCC Netgraph.

The distance between any two criminals in the undirected Netgraph ranges from 2-4 (see Table 4) while in the undirected BCC Netgraph it ranges from 1-7 (see Table 5). Using the shortest paths' lengths count,

there are variations in the connections formed between criminals in the undirected BCC Netgraph compared to the undirected Netgraph. In the undirected BCC Netgraph, there are some direct links between certain criminals, for example Andrew Fastow (BCCNet. ID 686) to Lea Fastow (11010), Andrew Fastow (686) to Ben Glisan (1369) and from Michael Kopper (12708) to Ben Glisan (1369). The emphasis of this paper is on the BCC Netgraph. In the next subsection, we describe our shortest paths network search algorithm briefly.

2.4. Criminal network formation methods

Past research shows that a criminal community can be formed using certain pre-defined rules. For example in (Al-Zaidy et al., 2012), a set of people belong to the same community if their names appear together in a document while (Anwar and Abulaish, 2014) group people into a community based on overlapping interests across different chat sessions. Our shortest paths network search algorithm is used to form a relationship network between suspects as a basis for an investigation (Magalingam et al., 2014). Unlike Al-Zaidy et al. (2012) and Anwar and Abulaish (2014), our algorithm does not restrict community membership to those with similarity in email content. By retaining duplicate email addresses (unlike Al-Zaidy et al. (2012)), we show that these could indicate secret trusted connections. Our network link doesn't depend on overlapping interest such as in Anwar and Abulaish (2014) but depends on a node's associations with particular central nodes and its links to known suspects. In our algorithm, the node or edge with highest centrality value is used. Girvan and Newman (Girvan and Newman, 2002; Newman and Girvan, 2004) remove the edge with highest betweenness score till they find different sub-networks or communities. In more recent work, (Ferrara et al., 2014) use a log analysis tool that adapts several community detection algorithms as well as utilising modular optimisation as done in Girvan and Newman algorithm (Girvan and Newman, 2002) and Newman's fast algorithm (Newman, 2004) to detect criminal organisations using phone call networks. Unlike the log analysis tool of (Ferrara et al., 2014) that implements Girvan and Newman's algorithm, we retain the central nodes, using them to form sub-networks of interest.

The shortest paths network search algorithm (SPNSA) is described in detail in (Magalingam et al., 2014). Firstly the algorithm requires a 'feed', a number of nodes of interest, whether these are known to be suspected criminals or otherwise regarded as relevant or important. For the Enron email network, each email account is used as a feed. For example, if a criminal has two email accounts, both the accounts are used in the feed list to represent that one criminal (see Table 1). Then the algorithm works by isolating a particular ego network, and within that ego network, identifying two central nodes, one with highest betweenness centrality and the other with highest eigenvector centrality. These nodes are named the Middle Man (MM) and the Most Influential (MI) respectively. Within the same ego network, the algorithm proceeds to extract the shortest paths from the ego to the central nodes, as well as from the ego to other nodes of interest in the list and finally, from the other nodes of interest to the central nodes. The three steps are repeated for every node of interest by selecting each in turn as the ego. The results of the various extractions are combined to

form a new sub-network. In Magalingam (2014), we studied the use of SPNSA on the directed BCC email network containing emails with 1 and 2 recipients bcc-ed. We found that the sub-network formed using SPNSA suggested possible people to investigate between known criminals and financial managers. This new sub-network could be used by an investigator as a preliminary investigative network (Magalingam et al., 2014). In this paper, we apply the SPNSA to analyse larger subsets of the Enron dataset than those studied in (Magalingam et al., 2014) in addition to comparing its performance against known community detection algorithms. Before applying the SPNSA, we use the different community detection algorithms to discover the various criminals’ communities in the Netgraph and BCC Netgraph.

3. Discovering criminals’ community using community detection algorithms

Many authors have used algorithms to analyse community structure and to consequently identify groups or sub-networks. An example of this is the use of network modularity in community detection algorithms (Girvan and Newman, 2002; Newman and Girvan, 2004). Communities exist when a graph consists of sets of nodes in tightly knit groups joined together by weaker connections between these groups (Girvan and Newman, 2002; Newman, 2010). The link structure and node attributes are the common components used in community detection algorithms. Girvan and Newman (Girvan and Newman, 2002; Newman and Girvan, 2004) repeatedly calculate edge betweenness, each time removing the edge with highest betweenness score such that as the graph becomes disconnected, the components represent each of communities. Other link based community detection approaches can be found in (Radicchi et al., 2004).

A different way of identifying communities is by using the node based approach called agglomerative algorithms (Blondel et al., 2008). Pons and Latapy (Pons, 2006) introduce a random walk concept which picks nodes from a network based on a fixed distance between two nodes. EAGLE is a software algorithm created by Shen et al. (Shen et al., 2009) that follows certain steps to form a community. First, it adapts the maximal clique calculation introduced by Bron and Kerbosch (Bron and Kerbosch, 1973), then removes the subordinate maximal clique. The algorithm then calculates the similarity between each pair of nodes in the clique, merges them into a new community, finds the similarity of the new community by comparing with an already existing community, repeating the steps until only one community remains.

Fastgreedy community detection (Clauset et al., 2004) uses a modularity optimization algorithm by first computing the fraction of within-community edges in a network, then subtracting from it the expected fraction of edges in a randomized version of the same network with same degree distribution. A nonzero value above 0.3 is considered a good measurement for the density of links inside communities (Clauset et al., 2004) (Blondel et al., 2008). Walktrap community detection merges similar nodes that are obtained using short random walks into a group (Pons, 2006). The leading eigenvector community detection algorithm implements the modularity optimization algorithm. It computes the modularity matrix and the eigenvector

of the matrix. It then divides the community based on the positive or negative sign of the elements in the eigenvector. If the large elements have the same sign, then the network has no community structure (Newman, 2006).

The results produced by applying these different methods to the Enron dataset to identify criminals' communities are discussed next. The community detection methods used are k -Neighbourhood, Fastgreedy, Walktrap and Leading Eigenvector algorithms all of which are available in the R igraph tool. Due to the connections between criminals being more visible in the undirected graph compared to the directed graph (See Tables 2, 3, 4 and 5) and since the majority of the community algorithms can only be applied to undirected graphs, we use the undirected Netgraph and BCC Netgraph for this exercise. We will later compare (in Section 4.2) these results with those obtained by using our shortest paths network search algorithm.

3.1. k -Neighbourhood detection

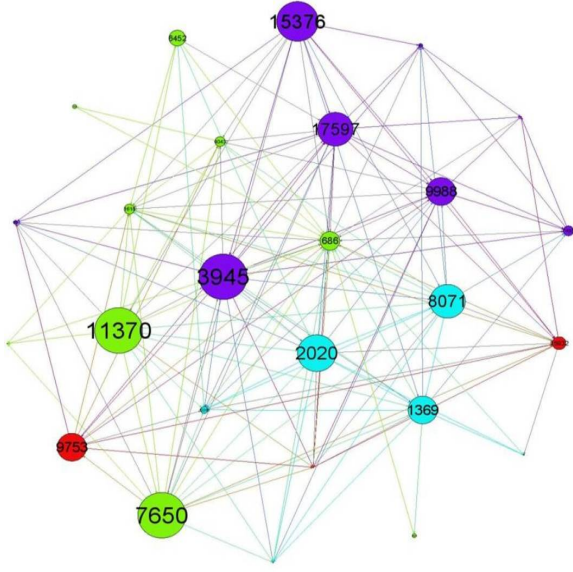
We first compute the total degree of each criminal in both the undirected Netgraph and BCC Netgraph. Table 6 shows the values. The neighbourhood function has been used previously to form a subgraph and identify the nearest link from a criminal node (Savage et al., 2014; Yasin et al., 2014). Using Table 6, we find all the neighbours of Andrew Fastow (686), (the first criminal in our list) in the undirected BCC Netgraph at a distance of 1 to 4 and form the networks.

Table 6: Enron money laundering criminals' degree

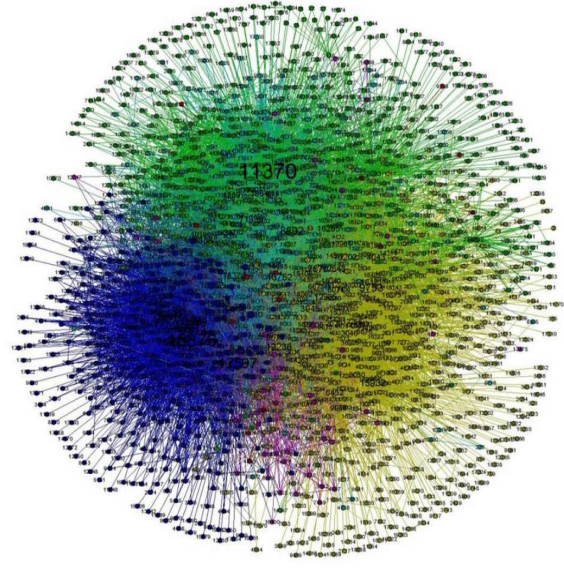
Name	Net. ID	Degree	BCCNet. ID	Degree
Andrew Fastow	1472	261	686	25
Andrew Fastow	-	-	687	1
Lea Fastow	17589	3	11010	4
Lea Fastow	17588	4	11009	4
Kevin Hannon	16202	1	10068	36
Kenneth Rice	16115	11	9994	4
Rex Shelby	23983	97	15224	21
Rex Shelby	23985	2	15225	2
A. Khan	-	-	205	2
Michael Kopper	20217	40	12708	6
Ben Glisan	-	-	1369	105
Joe Hirko	14052	11	8716	2
S. Yaeger	-	-	861	1

Table 6 shows the total degree of each criminal in Netgraph and BCC Netgraph.

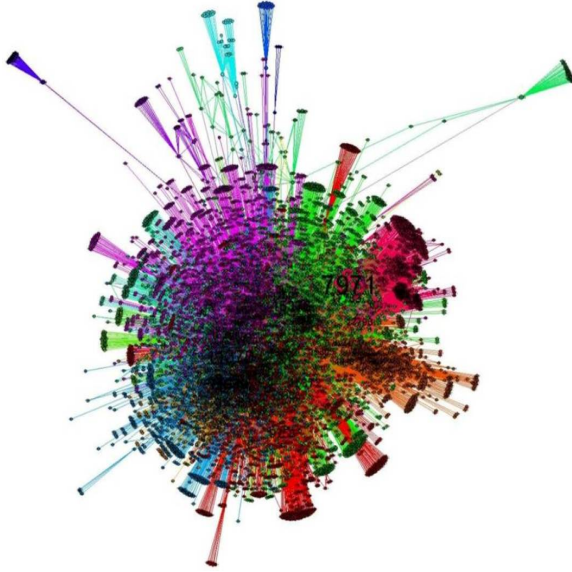
Figure 2 shows the networks found using 1,2,3,4-neighbourhood of Andrew Fastow (686) in the undirected BCC Netgraph respectively.



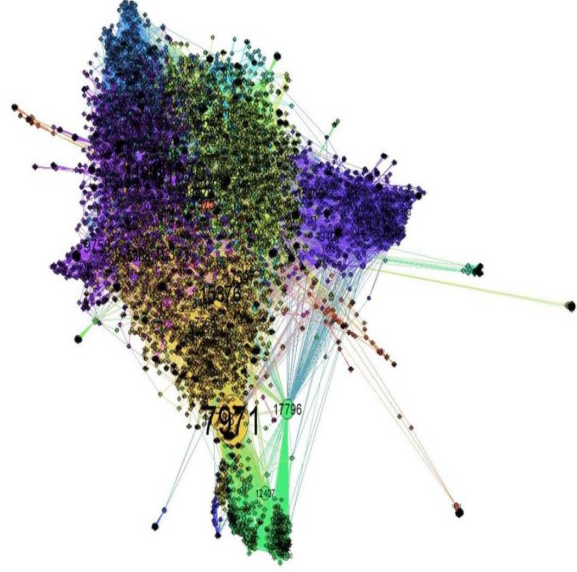
(a) 1-N network (26 nodes, 149 edges)



(b) 2-N network (1,559 nodes, 15,853 edges)



(c) 3-N network (8,433 nodes, 49,086 edges)



(d) 4-N network (14,916 nodes, 60,593 edges)

Figure 2: The figures above show the networks formed by using 1,2,3,4-neighbourhood (1,2,3,4-N) of Andrew Fastow (686) in the undirected BCC Netgraph respectively. In the 1-neighbourhood of Andrew Fastow only one other criminal was found, Ben Glisan (1369). The 2, 3, and 4- neighbourhood networks are clearly too dense to be able to identify other criminals easily.

In the 1-neighbourhood network of Andrew Fastow only one criminal was found, Ben Glisan (1369). According to Tables 2,3, 4,5 using either of the Andrew Fastow's email account (BCCNet. ID 686 or 687)

and either of the undirected or directed BCC Netgraphs, the one or two neighbourhoods would only rarely contain the other known criminals. The size of the network becomes bigger as the neighbourhood increases. The same method when applied to the undirected Netgraph, also produces large graphs that are difficult to explore. Clearly, using these dense network graphs, it is difficult to analyse a criminal’s occurrence and connections with other nodes.

3.2. Community detection algorithms in R igrph

In the next two sub-sections, 3.2.1 and 3.2.2, we compare the number of criminals, communities and connection of criminals to other nodes found using the community detection algorithms: Fastgreedy (Clauset et al., 2004), Walktrap (Pons, 2006) and Leading Eigenvector (Newman, 2006). As all three of these algorithms require undirected graphs, we use the undirected Netgraph and BCC Netgraph for this experiment.

3.2.1. Results of undirected Netgraph

Here we look at the results of applying the community detection algorithms to the undirected Netgraph. Applying the Fastgreedy community detection algorithm gives 15 communities. These communities range in size from 5,803 nodes to just 2 nodes. 6 out of 10 criminals were found in the second largest community that had 5,163 nodes and 22,936 edges. One criminal, Kevin Hannon (16202) appeared in a much smaller group consisting of 230 nodes and 781 edges.

Next, the Walktrap community detection algorithm was applied to the undirected Netgraph with the length of the random walk being 10 steps. 530 communities were found by the algorithm, with the first community detected consisting of 3,126 nodes and 12,462 edges, and again contained 6 of the 10 criminals (See Table 7). It was the second largest community formed by Walktrap. The largest community had 7,935 nodes while smallest one had just 1 node.

The third community detection algorithm used was the Leading Eigenvector. This detection algorithm detected 2 communities with the largest one containing 26,025 nodes and the smallest 2 nodes. All 7 criminals appeared in the largest community but the criminals were found to be isolated (See Table 7).

3.2.2. Results of undirected BCC Netgraph

The community detection algorithms were next applied to the undirected BCC network graph. The BCC Netgraph contains 65,532 edges and 19,716 nodes. The Fastgreedy algorithm found 832 communities, finding a number of small communities with less than 6 nodes each. 5 criminals were found to be in the largest community that had 2,195 nodes and 7,903 edges; Andrew Fastow (BCCNet. ID 686), Lea Fastow (BCCNet. ID 11010, BCCNet. ID 11009), Kevin Hannon (BCCNet. ID 10068), Ben Glisan (BCCNet. ID 1369) and Kenneth Rice (BCCNet. ID 9994). Ben Glisan (BCCNet. ID 1369) had the highest degree in this community. Meanwhile, Rex Shelby (15224, 15225) and S. Yaeger (861) belonged to the second largest

community with 2,142 nodes and 8,222 edges. The criminal who was in one of the smaller communities was Joe Hirko (8716).

Using the Walktrap community detection algorithm produced 1,773 communities from the undirected BCC Netgraph. The largest community had 1,493 nodes and the smallest one had just 1 node. Seven out of 10 criminals happened to exist in the same community that had 1,254 nodes (See Table 7). The Leading Eigenvector community detection algorithm found 719 communities in the BCC Netgraph. It ranged from the largest community with 15,792 nodes and the smallest with 1 node. Most of the criminals appeared in the largest community. In the Table 7, we list the communities where the criminals belong to and the community size that we identified manually.

Table 7: Criminals Found in Different Communities

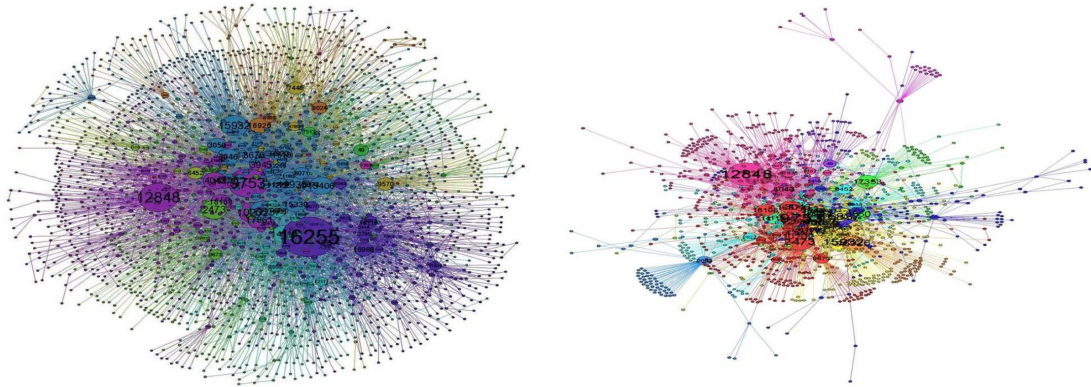
Net. ID	FG Com. ID	WT Com. ID	LEC Com. ID	BCCNet. ID	FG Com. ID	WT Com. ID	LEC Com. ID
1472	{7, 5163}	{1, 3126}	{1, 26025}	686	{5, 2195}	{36, 1254}	{1, 2001}
-	-	-	-	687	{5, 2195}	{594, 2}	{719, 15792}
17589	{7, 5163}	{1, 3126}	{1, 26025}	11010	{5, 2195}	{594, 2}	{719, 15792}
17588	{7, 5163}	{1, 3126}	{1, 26025}	11009	{5, 2195}	{36, 1254}	{719, 15792}
16202	{10, 230}	{29, 228}	{1, 26025}	10068	{5, 2195}	{36, 1254}	{719, 15792}
16115	{7, 5163}	{1, 3126}	{1, 26025}	9994	{5, 2195}	{36, 1254}	{719, 15792}
23983	{7, 5163}	{1, 3126}	{1, 26025}	15224	{3, 2142}	{36, 1254}	{719, 15792}
23985	{7, 5163}	{4, 7935}	{1, 26025}	15225	{3, 2142}	{1365, 1}	{719, 15792}
-	-	-	-	205	{224, 3}	{1030, 3}	{719, 15792}
20217	{7, 5163}	{1, 3126}	{1, 26025}	12708	{2, 2113}	{36, 1254}	{1, 2001}
-	-	-	-	1369	{5, 2195}	{45, 1493}	{719, 15792}
14052	{7, 5163}	{1, 3126}	{1, 26025}	8716	{17, 561}	{36, 1254}	{719, 15792}
-	-	-	-	861	{3, 2142}	{54, 1101}	{719, 15792}
Total community	15	530	2	-	832	1773	719

Table 7 shows community IDs to which each criminal belongs. The community and the size of each community is represented in curly brackets as $\{i_{th} \text{ community, size}\}$. The title of each column are Net. ID (Netgraph ID), FG Com. ID (Fastgreedy Community ID), WT Com. ID (Walktrap Community ID) and LEC Com. ID (Leading Eigenvector Community ID). The total number of communities formed is shown at the bottom of the table. The communities with smallest number of nodes, 1-3 nodes are highlighted in red.

3.2.3. Discussion of results obtained by R igraph community detection algorithms

The network partitioning using leading eigenvector to form communities seems not to be very effective for either the undirected Netgraph or BCC Netgraph, as the biggest community contains almost all the nodes. Some abnormal network structures were found in the communities of certain criminals (see Table 7). From the results of the Walktrap algorithm (highlighted in red), Andrew Fastow (687) and Lea Fastow (11010) belong to a small group of just two nodes, forming a community on their own. A. Khan (205) also belongs to a community of just three nodes. A. Khan is linked to two other nodes with email addresses; toriarules@aol.com and mmorales@arnel.com. These emails addresses are found to be external emails that do not belong to the Enron company email group. The Fastgreedy algorithm results in Rex Shelby (15225) being isolated in a community of his own, numbered 1365. All of these abnormalities occurred in the undirected BCC Netgraph. We also found more communities appearing in the undirected BCC Netgraph compared to the undirected Netgraph. The total number of communities formed using each detection algorithm is also shown in the last row of Table 7.

The community that contains the most criminals; community numbered 5 using Fastgreedy algorithm and community numbered 36 using Walktrap algorithm on the undirected BCC Netgraph were extracted. The detection using Walktrap algorithm on the BCC Netgraph yields the best result, some criminals appear in a small community on their own and 7 of 1254 members were criminals, but the result shows enormous number of nodes and links; networks (as shown in Figure 3) an investigator would need to analyse in order to find any connections between criminals and other nodes.



(a) Community numbered 5 using Fast Greedy algorithm (b) Community numbered 36 using WalkTrap algorithm

Figure 3: Figures above shows the community numbered 5 (2195 nodes) using fast greedy algorithm and community numbered 36 (1254 nodes) using walktrap algorithm on undirected BCC Netgraph.

Another method called clique percolation community detection developed by (Palla et al., 2005) was also used to identify the Enron criminal community. We found that the clustering coefficient values for both the

undirected Netgraph and BCC Netgraph were so low that the nodes did not converge to form communities. A high average clustering coefficient to a respective random network is needed to form sub-networks or clusters (Palla et al., 2005; Hills et al., 2009). In the next section, we apply our shortest paths network search algorithm to all four networks, the directed and undirected, Netgraph and BCC Netgraph.

4. Application of Shortest Path Network Search Algorithm

In (Magalingam et al., 2014), the shortest paths network search algorithm (SPNSA) was used to identify a trust network from a network of emails that have 1 and 2 bcc-ed recipients respectively to start an investigation. Here, we apply our SPNSA to the directed and undirected Netgraph and the BCC Netgraph. Different from (Magalingam et al., 2014), the directed and undirected BCC Netgraph used here contain all bcc-ed recipients. In this section, all the criminals in Table 6 are used as the feed for SPNSA. There are 7 criminals in the Netgraph and 10 criminals in BCC Netgraph (see Table 6).

4.1. Application of SPNSA on directed Netgraph and BCC Netgraph

Applying the SPNSA to the directed Netgraph results in all 7 criminals occurring in the extracted shortest paths network except for one email ID of Lea Fastow (Net. ID 17589) (see Figure 4). Lea Fastow's Net. ID 17589 that represents her second email address didn't appear in the sub-network because the node doesn't have an out-component that builds paths to other criminals or the central nodes.

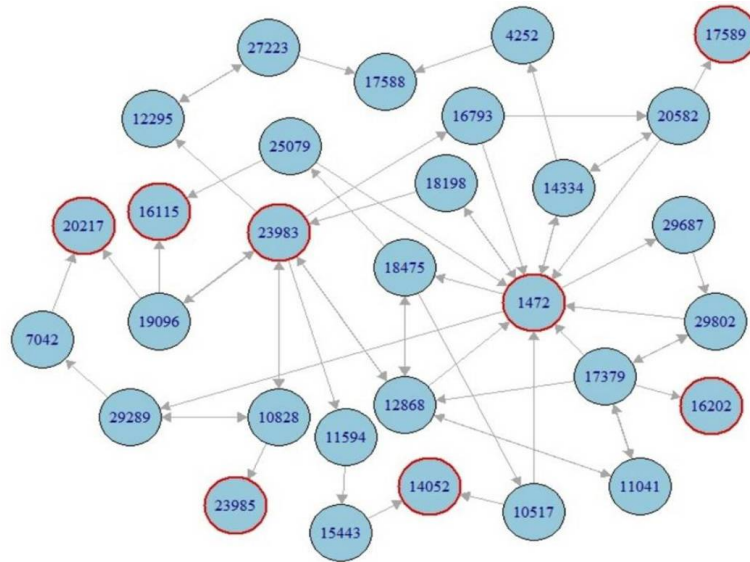


Figure 4: The shortest paths network formed using the directed Netgraph. All 7 criminals were found. The network formed is sparse and criminals' connections can be easily identified. This sub-network contains 30 nodes. The nodes highlighted in red represent the criminals with Net. ID in Table 1.

The sub-network of the directed BCC Netgraph captured using SPNSA is shown in Figure 5. This network contains 8 out of the 10 criminals. The two other criminals did not have any connection to other criminals or to the MM or MI, thus did not occur in this sub-network.

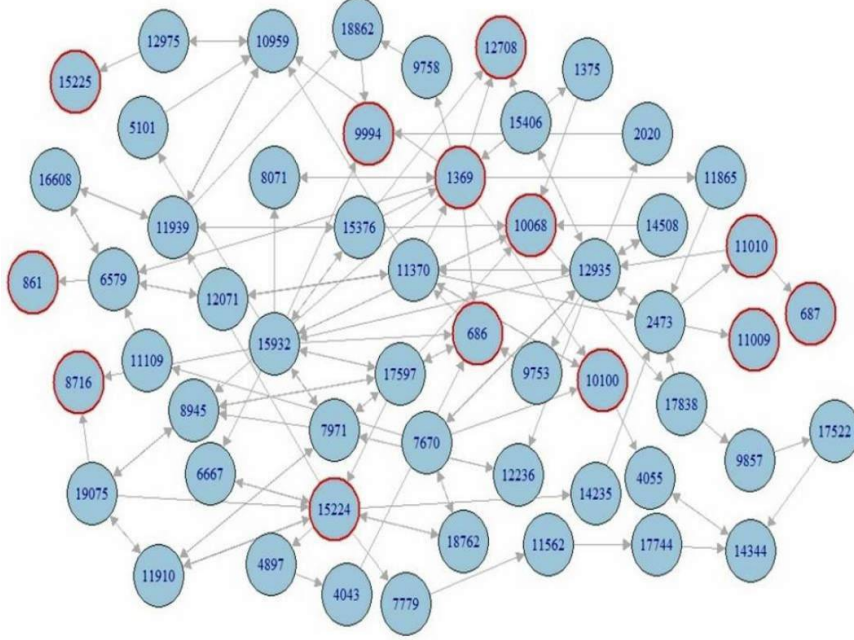


Figure 5: The shortest paths network formed using directed BCC Netgraph. 2 out of 10 criminals were lost. The network formed is sparse and criminals' connections can be identified. The number of nodes in this sub-network is 55. The nodes highlighted in red represent the criminals with BCCNet. ID in Table 1.

4.2. Application of SPNSA on the undirected Netgraph and the undirected BCC Netgraph

The SPNSA was then applied to the undirected Netgraph and the result is depicted in Figure 6. Again, all seven criminals occurred in this graph. This time all double email Net. IDs are captured in this sub-network.

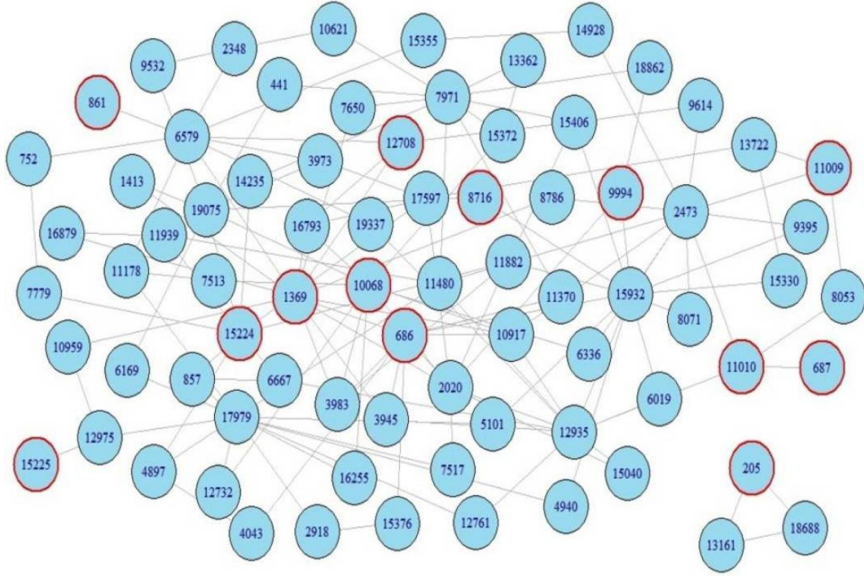


Figure 7: The shortest paths network formed using undirected BCC Netgraph. No criminal were lost. The network formed is sparse and criminals' connections can be identified. The size of this sub-network is 74 nodes and another small component of 3 nodes. The nodes highlighted in red represent the criminals with BCCNet. ID in Table 1.

The criminals and their links can be clearly seen in the results obtained (see Figures 4, 5, 6 and 7). The undirected BCC Netgraph yields the most number of criminals in the shortest paths network and connected components. A comparison of the results of using R *igraph* community detection algorithms with the result of shortest paths network search algorithm (SPNSA) applied to the undirected BCC Netgraph (see Figure 7) is documented in Table 8.

Table 8: Comparison between results found using R igraph community detection algorithms and SPNSA

	Community Detection Algorithm		Shortest Paths Network Search Algorithm	
	Undirected Netgraph	Undirected BCC Netgraph	Undirected Netgraph	Undirected BCC Netgraph
Community distribution	Nature: large and difficult to explore. Community Size: $250 < \text{nodes} < 26,100$.	Nature: large and difficult to explore. Community Size: $1,000 < \text{nodes} < 16,000$.	Nature: sparse to investigate and explore. Community Size: 30 nodes.	Nature: sparse to investigate and explore. Two components occurred, one of size 74 nodes and the other has 3 nodes.
Abnormalities	Andrew Fastow's external email add. (687) doesn't exist.	Walktrap community detection: found Andrew Fastow's external email add. (687) and Lea Fastow (11010) appeared in the same small community of size 2 nodes.	Andrew Fastow's external email add. (687) doesn't exist.	Found a direct connection between Andrew Fastow (687) and Lea Fastow (11010) that emerged in a community of size 74 nodes.
	A. Khan (205) doesn't exist	Fastgreedy and Walktrap community detection: found A. Khan (205) belongs to a small community of size 3.	A. Khan (205) doesn't exist.	Found A. Khan (205) belong to a small isolated community of size 3.
Total criminals	Detection Algorithm: Fastgreedy - 6/10 criminals in {7, 5163}. Walktrap - 6/10 criminals in {1, 3126}.	Detection Algorithm: Fastgreedy: 5/10 criminals in {5, 2195}. Walktrap: 7/10 criminals in {36, 1254}.	8/10 criminals.	All 10 criminals.

Table 8 shows the comparison between results found using R igraph community detection algorithms and SPNSA. The community formed by SPNSA is small and suitable for investigation. The 3-clique connection is formed as a separate network component and the nodes that are connected to the criminal can be easily identified. In the row giving total criminals, ($\{7, 5163\}$) refers to $\{i_{th} \text{ community, size}\}$ and the same follows for others.

5. Crime investigation methods using SPNSA

Anwar and Abulaish (2012) analysed their algorithm's performance by implementing three different scenarios based on the availability of information. Similar to (Anwar and Abulaish, 2012), here we specify certain ways an investigator could extract criminal subgraphs using the shortest paths network search algorithm for a preliminary investigation. In this section, the undirected BCC Netgraph is chosen instead of the undirected Netgraph due to three reasons found through our experiments in section 2.3 and the comparison between results in Table 8; more connections were detected between the criminals in the undirected BCC Netgraph (see Tables 2, 3, 4 and 5), cliques of two or three criminals were found in the undirected BCC Netgraph (See Table 8) and the most number of criminals were found in the undirected BCC Netgraph (See Table 8).

5.1. Extracting sub-networks using Non-Criminals

At an initial stage of a criminal investigation, an investigator may or may not have all or any of the criminals' details. The investigator could start the investigation with a suitable group of people. The

investigator will be able to feed in as many as necessary of the Enron managers' or suspects' node IDs in to the algorithm to form a network for their investigation. We simulate such a scenario by feeding in to the algorithm all the top managers in the Enron to obtain their communication network from the undirected BCC Netgraph. The result is shown in Figure 8.

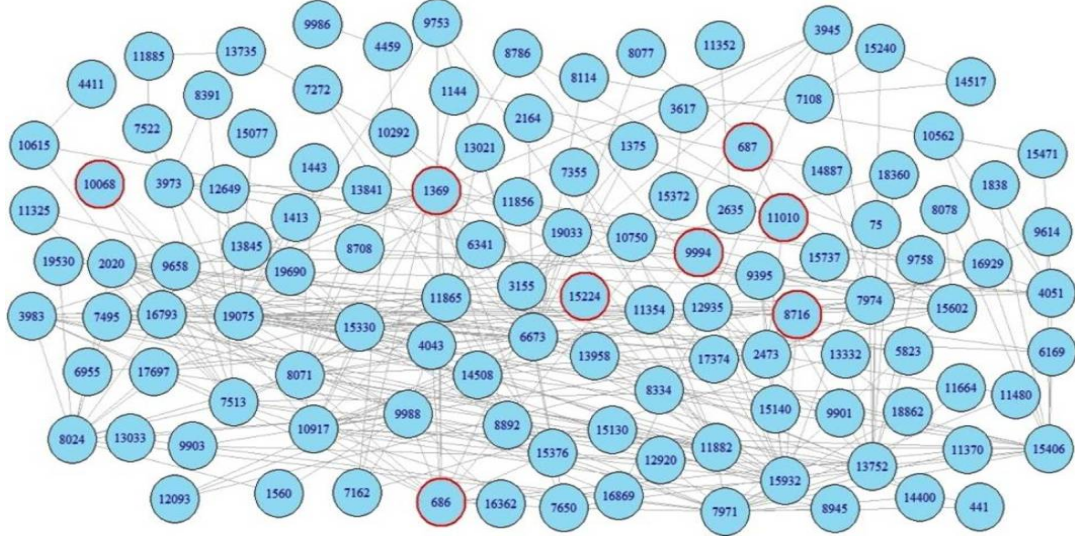


Figure 8: Enron undirected BCC Netgraph shortest paths network using all top managers as algorithm feed. The nodes highlighted in red are all criminals.

When compared to the network graph obtained in Figure 7, 7 out of 10 criminals are found by the algorithm this time around. Apart from the managers known to be criminals (see Section 2.2), the SPNSA also extracts 4 other criminals; Lea Fastow (11010), Kevin Hannon (10068), Rex Shelby (15224) and Joe Hirko (8716) with this top manager feed test.

Next a shortest paths network is formed using only financial managers. The financial managers' group is a subset of the top managers. The financial managers are the Head of Enron Global Finance, Sherron Watkins (16929), the Enron Chief Financial Officer, Andrew Fastow (686 , 687), the Enron Corporation Treasurer, Ben Glisan (1369), the Chief Accounting Officer, Rick Causey (15077), the Chief Financial Officer of Enron after Andrew Fastow, Jeff McMahon (8071). The network formed is shown in Figure 9.

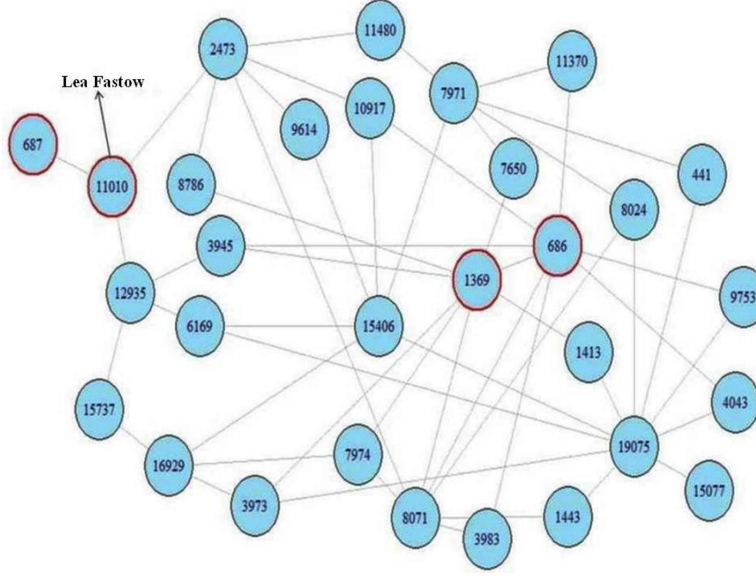


Figure 9: Enron undirected BCC Netgraph shortest paths network using all financial managers as algorithm feed. The nodes (686, 687 and 1369) highlighted in red are financial managers who are also criminals. The other criminal who was found is Lea Fastow (11010).

When comparing the criminals found in the sub-network formed using the Finance managers (see Figure 9), with those found in Figure 7, we see that Lea Fastow (11010) is the only other criminal found here, other than the finance managers who were also criminals.

5.2. Extracting subgraphs using leave-one-out method

The leave-one-out method is widely used in various fields of research as a data sampling method for an algorithm (Cawley and Talbot, 2004; Kocaguneli and Menzies, 2013) and can be used to estimate performance of a predictive model (Kocaguneli and Menzies, 2013). Past research (Shao, 1996) shows that one can set the number of data points to be removed from sample data and use it for validation. This is also called delete-p cross validation (Zhang, 1993).

We name this method as leave- C_i -out. Leave- C_i -out refers to dropping one criminal (C_i) from the list of criminals and running the shortest paths network search algorithm on the remaining criminals in the undirected BCC Netgraph. This method is a test of the ability of the algorithm to produce sub-networks that contain the convicted criminal not included in the algorithm feed. We name the criminal that is left out during each iteration as C_i . A criminal who has two email accounts has two different BCCNet. IDs and during the leave- C_i -out process both their IDs are dropped from the algorithm feed.

The results of the leave- C_i -out method are given in Table 9. In 5 out of 9 cases the criminal who is left out occurs in the network formed by the shortest paths algorithm. For example, we leave out Michael

Kopper (12708) from the feed (list of criminals) and run the algorithm. The result obtained is a sub-network that contains Michael Kopper and the connections of Michael Kopper (see Figure 10).

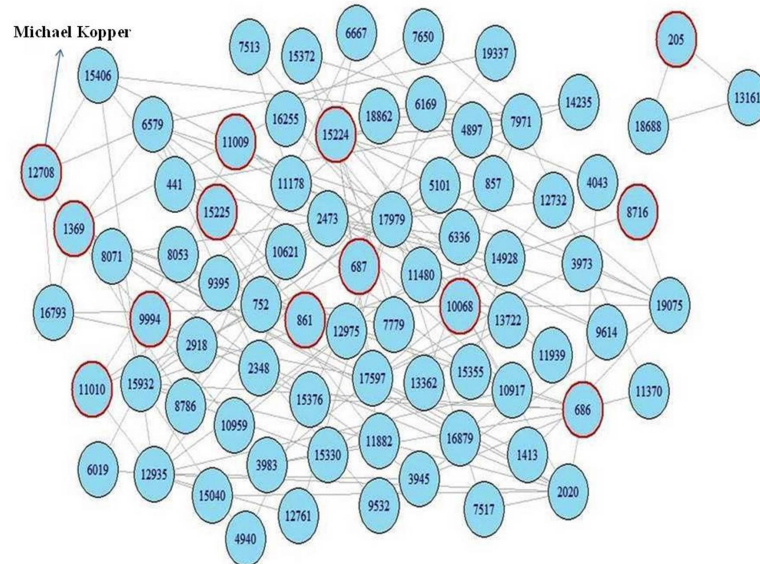


Figure 10: Enron undirected BCC Netgraph shortest paths network when Michael Kopper (12708) is left out from the criminal feed list. The nodes highlighted in red are all criminals.

Table 9: Leave- C_i -out Method

C_i Out	BCCNet. ID	C_i occur
Andrew Fastow	686, 687	✗
Lea Fastow	11010, 11009	✓
Kevin Hannon	10068	✗
Kenneth Rice	9994	✗
Rex Shelby	15224, 15225	✓
Michael Kopper	12708	✓
Ben Glisan	1369	✓
Joe Hirko	8716	✗
S. Yaeger	861	✓

Table 9 shows the result for Leave- C_i -out method. It is a test to see if the criminals that we left out still occur in the network formed by the shortest paths algorithm. C_i represents each criminal. The sign ‘✓’ indicates the C_i appeared in the shortest paths network community while ‘✗’ are given to C_i who do not appear in the extracted network. A. Khan (BCCNet. ID 205) was not included in this table, because, as shown in Figure 7, A. Khan (205) had no connections to other criminals, and further, appeared in a separate component containing 3 nodes.

6. Conclusion

The work presented in this paper contains the efficacy of the implementation of our shortest paths network search algorithm on larger dataset and the searches. The existing community detection algorithms in igraph did show the number of communities but an investigator would need to manually check the community to which a criminal belongs. Retrieving the neighbourhood sub-networks of the criminals in the community, as identified by the existing community detection algorithms, resulted in dense networks, which were hard to visualise and possibly, even harder to analyse. The criminals’ connections in these sub-networks were hard to view.

Our shortest paths network search algorithm (SPNSA) clearly shows the criminals’ connections to other nodes in all the sub-networks it extracted. Three different investigation methods were tested using the SPNSA; when the investigator knows all the criminals, when the investigator fails to detect one of the criminals and when the investigator is at the starting stage and doesn’t have any information about the criminals. In all three scenarios, the sub-network formed by SPNSA were sparse and hence, suitable for an investigator to see the connections as well as conduct further investigations. The SPNSA algorithm was able to extract and show the abnormalities through the sub-networks formed; components that contains criminals’ connections with other nodes and the 3-clique component can be easily detected.

The SPNSA allows the investigator to feed in the early suspects or suspicious entity into the criminal list of the algorithm, a function that is not available through other community detection algorithms. The quality of a criminal investigation can be improved when we can specify some inputs as in this algorithm; SPNSA could be a very useful preliminary investigation tool for an investigator.

References

References

- Al-Zaidy R, Fung B, Youssef AM, Fortin F. Mining criminal networks from unstructured text documents. *Digital Investigation* 2012;8(3):147–60.
- Anwar T, Abulaish M. Identifying cliques in dark web forums - an agglomerative clustering approach. In: *Intelligence and Security Informatics (ISI)*, 2012 IEEE International Conference on. 2012. p. 171–3. doi:10.1109/ISI.2012.6284289.
- Anwar T, Abulaish M. A social graph based text mining framework for chat log investigation. *Digital Investigation* 2014;11:349–62.
- Basu A. Social Network Analysis: A Methodology for Studying Terrorism. In: *Social Networking*. Springer; 2014. p. 215–42.
- Blondel VD, Guillaume JL, Lambiotte R, Lefebvre E. Fast Unfolding of Communities in Large Networks. *Journal of Statistical Mechanics: Theory and Experiment* 2008;2008(10):P10008.
- Bogawar PS, Bhoyar KK. Email mining: a review. *IJCSI International Journal of Computer Science Issues* 2012;9(1):429–34.
- Brickey KF. From Enron to WorldCom and beyond: Life and Crime After Sarbanes-Oxley. *Washington University Law Quarterly* 2003;81:357–402.
- Bron C, Kerbosch J. Algorithm 457: Finding all cliques of an undirected graph. *Commun ACM* 1973;16(9):575–7. doi:10.1145/362342.362367.
- Cawley GC, Talbot NL. Fast Exact leave-one-out cross-validation of Sparse least-squares Support Vector Machines. *Neural Networks* 2004;17(10):1467–75.
- Chen H, Chung W, Xu JJ, Wang G, Qin Y, Chau M. Crime Data Mining: A General Framework and Some Examples. *Computer* 2004;37(4):50–6.
- Christin N, Yanagihara SS, Kamataki K. Dissecting one click frauds. In: *Proceedings of the 17th ACM Conference on Computer and Communications Security*. New York, NY, USA: ACM; CCS '10; 2010. p. 15–26. doi:10.1145/1866307.1866310.
- Clauset A, Newman MEJ, Moore C. Finding community structure in very large networks. *Phys Rev E* 2004;70:066111. doi:10.1103/PhysRevE.70.066111.
- Cohen WW. Enron Email Dataset. 2009. URL: <http://www.cs.cmu.edu/~enron/>.
- Didimo W, Liotta G, Montecchiani F, Palladino P. An advanced network visualization system for financial crime detection. In: *Pacific Visualization Symposium (PacificVis)*, 2011 IEEE. 2011. p. 203–10. doi:10.1109/PACIFICVIS.2011.5742391.
- Ferrara E, De Meo P, Catanese S, Fiumara G. Detecting Criminal Organizations in Mobile Phone Networks. *Expert Systems with Applications* 2014;41(13):5733–50.
- Fox GS, Schaefer BE. Trusted electronic communications. 2012. US Patent App. 13/530,713.
- Girvan M, Newman MEJ. Community structure in social and biological networks. *Proceedings of the National Academy of Sciences of the United States of America* 2002;99(12):7821–6. doi:10.1073/pnas.122653799.
- Hills TT, Maouene M, Maouene J, Sheya A, Smith L. Categorical structure among shared features in networks of early-learned nouns. *Cognition* 2009;112(3):381–96.
- Iqbal F, Fung BC, Debbabi M. Mining Criminal Networks from Chat Log. In: *Web Intelligence and Intelligent Agent Technology (WI-IAT)*, 2012 IEEE/WIC/ACM International Conferences on. IEEE; volume 1; 2012. p. 332–7.

- Keila P, Skillicorn DB. Detecting unusual email communication. In: Proceedings of the 2005 conference of the Centre for Advanced Studies on Collaborative research. IBM Press; 2005. p. 117–25.
- Kocaguneli E, Menzies T. Software Effort Models should be Assessed via leave-one-out Validation. *Journal of Systems and Software* 2013;86(7):1879–90.
- Krebs VE. Uncloaking terrorist networks. *First Monday* 2002;7(4). URL: http://firstmonday.org/issues/issue7_4/krebs/index.html.
- Louis A, Engelbrecht AP. Unsupervised Discovery of Relations for Analysis of Textual Data. *Digital Investigation* 2011;7(3):154–71.
- Magalingam P, Rao A, Davis S. Identifying a Criminal’s Network of Trust. In: Third Int’l Workshop on Complex Networks and their Applications. Proceedings of the Tenth International Conference on Signal-Image Technology and Internet-Based Systems. 2014. p. 309–16.
- McDowell M, Householder A. Benefits of bcc. 2009. URL: <https://www.us-cert.gov/ncas/tips/ST04-008>.
- Murynets I, Piqueras Jover R. Crime scene investigation: Sms spam data analysis. In: Proceedings of the 2012 ACM conference on Internet measurement conference. ACM; 2012. p. 441–52.
- Nadji Y, Antonakakis M, Perdisci R, Lee W. Connected Colors: Unveiling the Structure of Criminal Networks. In: Research in Attacks, Intrusions, and Defenses. Springer; 2013. p. 390–410.
- Newman M. In: *Networks: An Introduction*. New York, NY: Oxford University Press; 2010. .
- Newman MEJ. Fast algorithm for detecting community structure in networks. *Phys Rev E* 2004;69:066133. doi:10.1103/PhysRevE.69.066133.
- Newman MEJ. Finding community structure in networks using the eigenvectors of matrices. *Phys Rev E* 2006;74:036104. doi:10.1103/PhysRevE.74.036104.
- Newman MEJ, Girvan M. Finding and evaluating community structure in networks. *Phys Rev E* 2004;69(2):026113. doi:10.1103/PhysRevE.69.026113.
- Oatley G, Crick T. Measuring uk crime gangs. In: Advances in Social Networks Analysis and Mining (ASONAM), 2014 IEEE/ACM International Conference on. 2014. p. 253–6. doi:10.1109/ASONAM.2014.6921592.
- Palla G, Derényi I, Farkas I, Vicsek T. Uncovering the overlapping community structure of complex networks in nature and society. *Nature* 2005;435(7043):814–8.
- Pons Pascal LM. Computing communities in large networks using random walks. *Journal of Graph Algorithms and Applications* 2006;10(2):191–218.
- Radicchi F, Castellano C, Cecconi F, Loreto V, Parisi D. Defining and identifying communities in networks. *Proceedings of the National Academy of Sciences* 2004;101(9):2658.
- Salter MS. *Innovation Corrupted: The Origins and Legacy of Enron’s Collapse*. Harvard University Press, 2008.
- Savage D, Zhang X, Yu X, Chou P, Wang Q. Anomaly Detection in Online Social Networks. *Social Networks* 2014;39:62–70.
- Securities U, Release ECP. Andrew S. Fastow, former Enron Chief Financial Officer, pleads guilty, settles civil fraud charges and agrees to cooperate with ongoing investigation. 2004. URL: <http://www.sec.gov/litigation/complaints/comp17762.htm>.
- Shao J. Bootstrap Model Selection. *Journal of the American Statistical Association* 1996;91(434):655–65.
- Shen H, Cheng X, Cai K, Hu MB. Detect overlapping and hierarchical community structure in networks. *Physica A: Statistical Mechanics and its Applications* 2009;388(8):1706 –12. doi:<http://dx.doi.org/10.1016/j.physa.2008.12.021>.
- Yang CC, Ng TD. Terrorism and Crime Related Weblog Social Network: Link, Content Analysis and Information Visualization. In: *Intelligence and Security Informatics, 2007 IEEE*. 2007. p. 55–8.
- Yasin M, Qureshi JA, Kausar F, Kim J, Seo J. A Granular Approach for user-centric Network Analysis to Identify Digital Evidence. *Peer-to-Peer Networking and Applications* 2014;:1–14.
- Zhang P. Model Selection via Multifold Cross Validation. *The Annals of Statistics* 1993;:299–313.

Zhou B, Pei J. The k-anonymity and l-diversity Approaches for Privacy Preservation in Social Networks against Neighborhood Attacks. Knowledge and Information Systems 2011;28(1):47–77.